

# Interview Questions

## Question 1: Explain the difference between agglomerative and divisive hierarchical clustering.

### Solution:

Agglomerative and divisive are two approaches to hierarchical clustering, which is a method used to group similar objects into clusters.

- **Agglomerative Clustering:** This is a bottom-up approach where each observation starts as its own cluster, and pairs of clusters are merged as one moves up the hierarchy. The process starts with  $n$  individual clusters (where  $n$  is the number of observations) and iteratively merges them based on a linkage criterion until all clusters have been merged into a single cluster that contains all observations. The key aspect of agglomerative clustering is the iterative process of finding and merging the closest (or most similar) pair of clusters until the desired number of clusters is reached or there is only one cluster left.
- **Divisive Clustering:** In contrast, divisive clustering is a top-down approach where all observations start in one cluster, and this cluster is recursively split into smaller clusters. The process begins with a single cluster containing all  $n$  observations, and at each step, a cluster is divided into smaller clusters based on a criterion that identifies the most dissimilar observations or groups within the cluster. This process continues recursively until each observation forms its own cluster or until a stopping criterion is met.

The main difference between the two methods is their direction of progress: agglomerative clustering builds clusters from the bottom up, while divisive clustering breaks down clusters from the top down.

## Question 2: What is a proximity matrix, and how is it used in hierarchical clustering?

### Solution:

A proximity matrix is a square matrix that gives the distance or similarity between each pair of observations in a dataset. In the context of hierarchical clustering, it's used to determine how close or similar each observation is to every other observation.

- In **agglomerative clustering**, the proximity matrix is crucial during the initial stage where each observation is considered as a separate cluster. The matrix is used to identify the two closest clusters to merge at each step of the algorithm. After each merger, the proximity matrix is updated to reflect the distances between the newly formed cluster and all other clusters.
- In **divisive clustering**, although the initial approach starts with a single cluster, the proximity matrix can still play a role in determining how the cluster is divided, by identifying the most dissimilar observations or subgroups within the cluster that should be split.

The proximity matrix can be based on various distance measures, such as Euclidean distance, Manhattan distance, or any other metric that suits the nature of the data and the specific requirements of the analysis.

## Question 3: Describe the implications of choosing different linkage criteria in agglomerative clustering.

### Solution:

The choice of linkage criteria in agglomerative clustering significantly affects the resulting cluster structure. Linkage criteria determine how the distance between clusters is calculated, and hence, which clusters are merged at each step. The most common linkage criteria are:

- **Single Linkage (Nearest Point):** The distance between two clusters is defined as the shortest distance from any member of one cluster to any member of the other cluster. This can lead to a "chaining effect" where clusters may end up being long and straggly.
- **Complete Linkage (Farthest Point):** The distance between two clusters is the longest distance from any member of one cluster to any member of the other cluster. This tends to produce more compact and well-separated clusters.
- **Average Linkage:** The distance between two clusters is the average distance between all pairs of members in the two clusters. This offers a balance between the single and complete linkage methods and can be less susceptible to outliers than the single linkage.
- **Ward's Method:** The distance between two clusters is the increase in the total within-cluster variance after merging the two clusters. This method tends to produce clusters of similar sizes and is particularly useful when the clusters have a roughly spherical shape.

The choice of linkage criterion affects the shape and size of the clusters formed. For instance, single linkage may be preferable for capturing non-elliptical shapes, while Ward's method is suitable for globular clusters. Therefore, the choice should be based on the specific characteristics of the data and the analytical goals.

#### **Question 4: How do you read a dendrogram, and what does it tell you about the data?**

##### **Solution:**

A dendrogram is a tree-like diagram that illustrates the arrangement of the clusters produced by hierarchical clustering. To read a dendrogram:

- **Vertical Lines:** Represent clusters being merged (in agglomerative) or split (in divisive). The height of the merge (or split) point indicates the distance (dissimilarity) at which the clusters were combined (or separated).
- **Horizontal Lines:** Connect clusters at the level of their similarity or dissimilarity. The length of the horizontal lines does not have a direct interpretation but serves to connect vertical lines.
- **Reading Clusters:** To determine the number of clusters at a given level of similarity, draw a horizontal line across the dendrogram. The number of vertical lines it intersects gives the number of clusters at that similarity level.

Dendrograms provide valuable insights into the data, such as the hierarchical relationship between clusters, the relative similarity between observations, and potential groupings within the data. By analyzing the structure of the dendrogram, one can also infer about the natural grouping in the data and decide on a reasonable number of clusters by identifying significant jumps in dissimilarity (represented by the height of the vertical lines).

#### **Question 5: Discuss the limitations of agglomerative clustering and how they might affect the analysis.**

##### **Solution:**

Agglomerative clustering, while widely used, has several limitations that could impact its effectiveness in certain analyses:

- **Scalability:** The computational complexity of agglomerative clustering can be quite high, especially for large datasets. This is because, at each step, the algorithm must update the distance matrix and identify the closest pair of clusters, which can be computationally intensive.
- **Sensitivity to Noise and Outliers:** Agglomerative clustering can be sensitive to noise and outliers because the presence of outliers can distort the distance calculations, leading to non-representative clusters.
- **Irreversibility:** Once two clusters are merged, the decision is final and cannot be undone. This can lead to suboptimal clustering if early merges are not optimal.
- **Choice of Linkage Criteria:** The results of agglomerative clustering can vary significantly based on the choice of linkage criteria. There's no one-size-fits-all criterion, and the choice can greatly affect the shape and size of the resulting clusters.

These limitations mean that while agglomerative clustering is a powerful tool for exploratory data analysis, careful consideration must be given to the nature of the data, the scale of the dataset, and the choice of linkage criteria to ensure that the results are meaningful and relevant to the analytical objectives.